

AD-A205 243

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

DTIC FILE COPY

(2)

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER AR6 22942.14-MA	2. GOVT ACCESSION NO. N/A	3. RECIPIENT'S CATALOG NUMBER N/A
4. TITLE (and Subtitle) Investigations in Robust/Resistant Data Analysis: Final Report		5. TYPE OF REPORT & PERIOD COVERED final report 1 Nov 85 - 31 Oct 88
7. AUTHOR(s) David C. Hoaglin and Frederick Mosteller		6. PERFORMING ORG. REPORT NUMBER AR-126
9. PERFORMING ORGANIZATION NAME AND ADDRESS Department of Statistics, Harvard University One Oxford Street Cambridge, MA 02138		8. CONTRACT OR GRANT NUMBER(s) DAAG29-85-K-0262
11. CONTROLLING OFFICE NAME AND ADDRESS U. S. Army Research Office Post Office Box 12211 Research Triangle Park, NC 27709		10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)		12. REPORT DATE 27 December 1988
		13. NUMBER OF PAGES 18
		15. SECURITY CLASS. (of this report) Unclassified
		15a. DECLASSIFICATION/DOWNGRADING SCHEDULE
16. DISTRIBUTION STATEMENT (of this Report) Approved for public release; distribution unlimited.		
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report) NA		
18. SUPPLEMENTARY NOTES The view, opinions, and/or findings contained in this report are those of the author(s) and should not be construed as an official Department of the Army position, policy, or decision, unless so designated by other documentation.		
19. KEY WORDS (Continue on reverse side if necessary and identify by block number) analysis of variance, distribution shape, exploratory data analysis, interaction, outliers, regression diagnostics, robustness, statistical software		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number) Research under this contract involved three main problem areas: (1) diagnostic and critical data analysis (including exploratory and robust methods), (2) analysis of distribution shape (techniques and applications to real data), and (3) implementation in computer software of selected data-analytic techniques. This final report summarizes the research and major results in each of these three areas.		

DTIC
ELECTE
S FEB 15 1989 D
H

INVESTIGATIONS IN ROBUST/RESISTANT DATA ANALYSIS
FINAL REPORT

David C. Hoaglin
and
Frederick Mosteller

27 December 1988

U. S. ARMY RESEARCH OFFICE

Contract DAAG29-85-K-0262

with

Harvard University

APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNLIMITED.

89 2 15 058

THE VIEWS, OPINIONS, AND/OR FINDINGS CONTAINED IN THIS REPORT ARE THOSE OF THE AUTHORS AND SHOULD NOT BE CONSTRUED AS AN OFFICIAL DEPARTMENT OF THE ARMY POSITION, POLICY, OR DECISION, UNLESS SO DESIGNATED BY OTHER DOCUMENTATION.

Contents

Problems and Results	1
Diagnostic and Critical Data Analysis	1
Distribution Shape	7
Computer Software	10
Publications and Technical Reports	11
Participating Scientific Personnel	13
Bibliography	14



Accession For	
NTIS GRA&I	<input checked="checked" type="checkbox"/>
DTIC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	
By	
Distribution/	
Availability Codes	
Dist	Avail and/or Special
A-1	

PROBLEMS AND RESULTS

This contract continued ongoing research in three main problem areas: (1) diagnostic and critical data analysis, (2) distribution shapes in real data, and (3) computer implementation of selected data-analytic techniques.

Diagnostic and Critical Data Analysis

As data analysis proceeds in an exploratory mode, it often raises questions about the strength of patterns that have been uncovered or about individual observations that seem unusual. Critical data analysis aims to deal with these types of confirmatory questions.

During the period of the contract our research in this area involved outlier detection, robust analysis of variance, regression diagnostics, robust/resistant techniques in quality control, and other applications.

In earlier work related to outlier detection, Hoaglin, Iglewicz, and Tukey (1986) studied the performance of a standard outlier-labeling technique from exploratory data analysis. This rule uses the lower fourth F_L and the upper fourth F_U of the sample (approximate sample quartiles) as the basis for the cutoffs

$$F_L - k(F_U - \tilde{F}_L) \text{ and } F_U + k(F_U - F_L) \quad (1)$$

with $k = 1.5$. Any observations that fall outside these cutoffs are labeled as possible outliers and subjected to closer scrutiny whenever feasible.

This work measured the performance of this and other rules in terms of the *some-outside rate per sample*: the probability that a sample of n contains at least one observation beyond the cutoffs. By

design the basic exploratory rule ($k=1.5$) has a relatively high some-outside rate per sample for Gaussian data: roughly between 15 percent and 35 percent as n ranges from 5 to 50. The original definition gives the fourths in terms of the depth, $f = \frac{1}{2}[(n+3)/2]$ with $[\cdot]$ the greatest-integer function, and the sample order statistics $X_{(1)} \leq \dots \leq X_{(n)}$ according to

$$F_L = X_{(f)} \text{ and } F_U = X_{(n+1-f)}.$$

As a consequence the some-outside rate per sample increases as n increases, following separate but similar curves for the four values of $n \bmod 4$. This behavior may be undesirable if one does not insist on a simple definition of f that facilitates hand calculation.

Hoaglin and Iglewicz (1987b) studied the results of fine-tuning such outlier-labeling rules in two ways. First, they considered two forms of smooth interpolation for the fourths: the "ideal" definition $f_I = n/4 + (5/12)$ and an alternative based on a common definition of the quartiles, $f_Q = n/4 + (1/4)$. With these definitions the some-outside rate per sample follows a single smooth curve against n . Second, to permit use of such rules for more traditional outlier detection, they determined values of the constant k (for f_I , f_Q , and the standard definition, and at selected sample sizes $n \leq 300$) that maintain the some-outside rate per sample at .10 or .05 for Gaussian data. As a rough approximation the rule based on f_I can use $k = 2.2$ out to about $n = 50$ for a rate of .05. For a rate of .10 no constant value of k seems satisfactory; a crude approximation yields $k \approx 2.02 - 2.5/n$ for $7 \leq n \leq 52$.

In addition to their simple form, outlier-detection rules that use resistant cutoffs (as in equation (1)) have the convenient property that, within reasonable limits, the user does not need to specify the maximum possible number of outliers. The limitation arises from the breakdown point of the location measure and scale measure used in specifying the cutoffs. For a measure of location, say, breakdown involves replacing observations in the sample by arbitrarily extreme values. The

breakdown point is the smallest fraction of such contamination that can cause the value of the measure to become arbitrary. Because of the way that they use the fourths, the cutoffs in equation (1) have a breakdown point of essentially 25%.

One can also start with a robust/resistant location estimator T and a robust/resistant scale estimator S and then use these to define an outlier-labeling rule by placing cutoffs at $T \pm kS$ for some suitable positive constant k . By choosing particular T and S that have good efficiency as robust estimators, one hopes to produce an outlier-labeling rule (and corresponding outlier-detection rule) that performs well according to a variety of criteria. Hoaglin and Iglewicz (1987a) investigated rules based on a biweight M-estimator of location, T_{bi} , and a biweight A-estimator of scale, s_{bi} (see, for example, Kafadar (1983) and Lax (1985)). To facilitate comparison, they matched the population values of the cutoffs $T_{bi} \pm k's_{bi}$ in Gaussian data to those of the exploratory rules in equation (1) by taking $k' = 0.6745 + 1.349k$. The resulting some-outside rate per sample (estimated by simulation) for Gaussian data is about 11 percentage points lower for $n \leq 50$ than the rate for the standard exploratory rule and about 7 percentage points lower than that for the exploratory rule with f_1 as the depth of the fourths. That is, when the data are ideally well-behaved, the biweight rule less often labels any observations as outside. Hoaglin and Iglewicz also studied other characteristics of the biweight rules; and they determined values of k (again at selected $n \leq 300$) that correspond to a some-outside rate per sample of .10 and .05, so that the rule can be used for outlier detection.

Thus far all the results described pertain to null situations, in which the data come from a Gaussian distribution. Hoaglin, Iglewicz, and Tukey (1986) included, as alternative null situations, some symmetric distributions with heavier tails, but none of the data have contained known outliers. To pursue the performance of various rules in detecting specified outliers (one suitable definition of power in this context), Hoaglin and Iglewicz (1988) developed a class of fixed-outlier models. A model of this type begins with r fixed observations and completes a "sample" of n by taking $n-r$ random observations from the standard Gaussian distribution.

For $r=1$ and the fixed observation at x , one determines the probability that a rule labels the observation at x as an outlier. By systematically varying x one obtains the rule's performance curve, $P_n(x)$. When $r=2$, the performance curve generalizes to two performance surfaces. If x_1 and x_2 are the fixed observations, $P_{n1}(x_1, x_2)$ is the probability that the rule labels x_1 an outlier, and $P_{n2}(x_1, x_2)$ is defined similarly for x_2 . For some rules the task of studying the performance surfaces simplifies, because the cross-sections at $x_2 = x_1$, $x_2 = -x_1$, and $x_2 = 0$ appear to capture the important features of P_{n1} and (by symmetry) P_{n2} . Thus one plots $P_{n1}(x, x)$, $P_{n1}(x, -x)$, and $P_{n1}(x, 0)$ against x . This approach provides control over the location of the outliers; and it is convenient for simulation, because one can generate the random part of the sample and then vary x over a grid.

Hoaglin and Iglewicz used the models with 1 and 2 fixed outliers to compare the performance of several outlier-detection rules at $n=20$ and $n=40$. They considered primarily the exploratory rule, the biweight rule, and a standard test based on the sample kurtosis. (For Gaussian data the kurtosis test has certain optimality properties when the outliers arise from arbitrary shifts in location or from arbitrary increases in variance.) All three were set for outlier detection at the .05 level. At both $n=20$ and $n=40$, $P_n(x)$ for the biweight rule follows that for the kurtosis test fairly closely, whereas the curve for the exploratory rule drops well below for $x > 3$. For the cross-section $P_{n1}(x, 0)$ a similar pattern emerges, as one would expect from the fact that the observation at 0 cannot be an outlier. The other two cross-sections give a more complicated picture. When the two fixed outliers are both at x and $n=20$, the biweight rule and the exploratory rule perform as well as the kurtosis test and perhaps slightly better over $3 \leq x \leq 5$. At $n=40$, however, the kurtosis test outperforms the other two by a substantial margin over the same interval. And similar patterns of domination emerged with the two fixed outliers at x and $-x$, both at $n=20$ and $n=40$. These results demonstrate the usefulness of the fixed-outlier model in comparing outlier-detection rules. Further study of them may lead to improvements in the biweight rule or suggest other robust/resistant rules.

Rocke (1983) studied robust estimation of variance components in the analysis of variance.

Using a substantial simulation, he compared the performance of the standard estimators of variance components with those of two robust methods (Huber and biweight) in the two-way mixed model with various relative sizes of between-means variance, various sizes and frequency of outliers, and various sample sizes. Rocke observed that for small between-means variance the biweight estimator is strongly biased compared to the Huber estimator. From further analysis of the results of Rocke's simulations, Mosteller found that, in large parts of the space of simulations, the biweight was much to be preferred in spite of its poor performance in the more "null" situation. More important, in much of the simulation space both estimators performed poorly, even when one was much to be preferred. This latter finding is especially troubling because the basic methods underlying these two estimators have been especially successful in handling outliers in the one-sample problem. In correspondence with John W. Tukey these results led us to conclude that we will need to reformulate the goals of robust estimation for the analysis of variance.

Because of their continuing interest in regression diagnostics for high-leverage and influential observations, Hoaglin and Kempthorne (1986) contributed discussion on a review paper by Chatterjee and Hadi. The discussion emphasized cutoffs, rules of thumb, and their role in identifying influential observations; simple residual plots that display high-leverage, outlying, and influential observations simultaneously; approaches to uncovering influential groups of observations; and selection of subsets of carriers. It also sketched a step-by-step diagnostic strategy that should be useful in practice.

As a basis for better empirical understanding of the behavior of regression diagnostics in large data sets, Hoaglin obtained (from a colleague at Abt Associates Inc.) the values of various regression diagnostics (e.g., diagonal elements of the hat matrix, DFITS, DBETAS, and studentized residuals) that had been computed for two sizable regressions. One involved $n = 1034$ observations, the other had 5719, and both had $p=21$ explanatory variables (including the constant). For the smaller data set he analyzed the frequency with which the values of DFITS and $DBETAS_1, \dots, DBETAS_{21}$ fell outside the recommended cutoff values ($\pm 2\sqrt{p/n}$ for DFITS and $\pm 2/\sqrt{n}$ for DBETAS). For DFITS the

percentage was 6.0%, whereas for the DBETAS it ranged from 0.6% to 6.8% (often closer to 4% or 6% than to the nominal 5%). Normal probability plots for DFITS and one of the DBETAS revealed that their distributions (across the observations in the data set) had substantially heavier tails than the Gaussian.

Basic techniques for diagnosing leverage and influence in regression could be applied more often than they so far seem to be. As one step in the process of speeding the spread of diagnostics in everyday applications of regression, Hoaglin (1988a) gave a tutorial account of leverage and influence in simple linear regression.

The area of quality control offers many opportunities to apply exploratory and robust methods. Iglewicz and Hoaglin (1987) devised a control chart that uses boxplots to show more information than the customary means and ranges and to reveal possibly outlying observations.

In an application of robust/resistant methods to anthropometry, Himes and Hoaglin (1989) investigated techniques for smoothing reference data on human growth. For such measurements as triceps skinfold thickness, the available (cross-sectional) data usually give selected percentiles by sex and single year of age. By working with the sequences of differences between adjacent age-specific percentiles, a variant of the exploratory technique known as delineation ensures that the resulting smoothed percentiles maintain the proper order. In an example Himes and Hoaglin obtained smoothed percentiles that provided a more satisfactory representation of the raw percentiles than did the published smoothed percentiles, which researchers at the National Center for Health Statistics had produced by straightforward use of cubic-spline regression.

For data that might be handled by two-way analysis of variance with one observation per cell, diagnosis of multiplicative structure can encounter difficulties if not done resistantly. In its most general form a single multiplicative term in the residuals after a simple additive fit, $e_{ij} = y_{ij} - (m + a_i + b_j)$, can be summarized as kc_id_j . Least-squares fitting, however, allows a single aberrant

observation among the y_{ij} to perturb the residuals in a pattern of the form kc_id_j . Thus an isolated anomaly leaks into the e_{ij} as a systematic pattern. One can avoid such leakage by obtaining the residuals from median polish or another suitably resistant analysis. To facilitate diagnosis of multiplicative structure, Hoaglin refined this approach by using the graphical display known as the scatterplot matrix. One examines the scatterplots for all pairs of columns (or rows) of the e_{ij} ; that is, the (s,t) cell in the matrix contains the scatterplot of e_{is} versus e_{it} , one point for each i . If the e_{ij} resemble kc_id_j , then the (s,t) scatterplot will tend to follow a straight line with slope d_s/d_t . One can readily assess the possibility of leakage problems by comparing the scatterplot matrix for the least-squares residuals with one for resistant residuals. Several examples have given encouraging results.

Distribution Shape

In a number of contexts it can be beneficial to develop detailed information on the shapes of distributions that tend to arise in real data. For example, data in a particular area of application may suggest a need for a robust estimator, and one might like to base the choice of an estimator on an assessment of the ways in which the underlying distribution may depart from an ideal shape (usually the Gaussian). In another example the samples may be large enough to provide a reasonable amount of information on the shape of the distribution, so that one can criticize an assumed distributional model and adapt the model to the data.

As part of our research on distribution shape we added substantially to our collection of data for study. The data sets include the following:

Opthamology (spherical and cylindrical refractive errors, $n = 510$)

Blood pressure (by sex in three age groups, n ranges from 534 to 965)

Laboratory measurements (mostly blood chemistry) on subjects in two treatment groups ($n \approx 1300$)

Annual incomes reported by low-income households in two sites ($n = 994$ and 608)

Climatic variables in Texas (15 variables for each of 189 quadrats covering the state)

Birth weights from the Medical Birth Registry of Norway: first and second births to women who had exactly two births, both singleton, during the period 1975-1984 (by the sexes of the two children and the year of the first birth, $n \approx 4000$ per stratum).

For the last of the data sets in this list, the Norwegian birth data, Hoaglin began systematic analysis, focusing primarily on questions related to low-weight births. He used a variety of graphical techniques including quantile-quantile plots, along with numerical techniques based on the g-and-h distributions (see below). Initial results indicated, among other things, that distributions of second-birth weights closely resemble a mixture of a Gaussian distribution and a small percentage of lower-weight births, but not so closely as previous analyses in the literature have assumed. He also found that adjustment of the birth weight of the second child for multiple linear regression on the birth weight of the first child and the gestational age of the second child produced a distribution whose shape is much closer to Gaussian. These shapes showed good agreement across the four combinations of the sexes of the two children and across strata on the birth weight of the first child. The findings should contribute to the discussion of the phenomenon of repeated low-weight births.

Focusing on another aspect of distribution shape, Mosteller gathered a substantial number of references about the frequency with which outliers enter real sets of data in the empirical literature, as well as about fitting the tails of data distributions. In further work he plans to study tail behavior by empirical fitting of generalized Pareto distributions (Smith, 1987).

For quantitative description of distribution shape, we often use Tukey's family of g-and-h distributions, discussed by Hoaglin (1985). By using quantiles the basic techniques for this family provide resistance, and they offer greater flexibility than the usual third and fourth moments. In terms of a standard Gaussian random variable Z the basic random variable Y of the g-and-h family (for a fixed value of the skewness parameter g and the elongation parameter h) is given by

$$Y = \frac{e^{gZ} - 1}{g} e^{hZ^2/2} . \quad (2)$$

One can readily introduce location and scale parameters.

To provide a description of skewness, for example, one begins with the median $y_{.5}$ and some selected pairs of quantiles y_p and y_{1-p} at tail areas p ($0 < p < .5$) that are either specified in advance (e.g., $p = \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \dots$) or suggested by the data. The position of y_p and y_{1-p} relative to $y_{.5}$ corresponds to a value of g (denoted by g_p or \hat{g}_p) according to

$$g_p = -\frac{1}{z_p} \log_e \frac{y_{1-p} - y_{.5}}{y_{.5} - y_p} , \quad (3)$$

where z_p is the p th quantile of the standard Gaussian distribution.

As a basis for our empirical work on distribution shape, we needed an estimate of the variability of \hat{g}_p from a set of data. Thus Hoaglin and Tukey (1988) developed an approximation for $\text{var}(\hat{g}_p)$. Their approach applies the anglit transformation to estimate the variance of \hat{y}_p , \hat{y}_{1-p} , and $\hat{y}_{.5}$. The data may come either as order statistics or as a frequency distribution. Hoaglin used these approximations in studying skewness in the Norwegian birth weight data.

Analyses of data often involve a transformation, such as the logarithm, square root, or reciprocal, and these nonlinear transformations affect the shape of distributions. To illustrate the point that such transformations or their results arise in everyday life more often than most people realize, Hoaglin (1988b) collected and discussed a variety of nontechnical examples, including magnitudes of stars, the Richter scale for earthquakes, intensity of sounds, average speed in auto races, and measures of gasoline consumption.

Computer Software

In connection with our work on new techniques of data analysis, we often develop computer software that implements the techniques. From the start the design of such programs looks beyond an initial version for personal use to a more polished version that could readily be made available to others.

During the period of the contract our main effort in this area focused on programs that calculate estimates of $\text{var}(\hat{g}_p)$ at selected values of p for data to which one is applying the g-and-h distributions, as described above. Separate versions start with an ordered sample or a frequency distribution. Each delivers the values of the \hat{g}_p , the estimated $\text{var}(\hat{g}_p)$, confidence intervals for the \hat{g}_p , and estimated covariances of the \hat{g}_p at different p .

Occasionally our interest in software leads us to review available implementations of a particular technique. Frigge, Hoaglin, and Iglewicz (1987, 1989) examined the implementations (in Minitab, S, SAS, SPSS, Statgraphics, and Systat) of the exploratory display known as the boxplot, discussed the confusing variety of nonstandard definitions, and made suggestions for improvements.

PUBLICATIONS AND TECHNICAL REPORTS

(Reports submitted for publication have been omitted when superseded by a published version.)

D. C. Hoaglin, "Influence and Collinearity in Regression: What Students Need to Know about Regression Diagnostics." Report AR-95, 6 June 1986.

D. C. Hoaglin, "Poissonness Plots and Related Techniques for Studying Discrete Frequency Distributions" (Fortran software). Report AR-97, 27 June 1986.

P. Langenberg and B. Iglewicz, "Trimmed Mean \bar{X} and R Charts," Journal of Quality Technology, 18 (1986), 152-161.

D. C. Hoaglin and P. J. Kempthorne, Discussion of "Influential Observations, High Leverage Points, and Outliers in Linear Regression" by S. Chatterjee and A. S. Hadi, Statistical Science, 1 (1986), 408-412.

D. C. Hoaglin, B. Iglewicz, and J. W. Tukey, "Performance of Some Resistant Rules for Outlier Labeling," Journal of the American Statistical Association, 81 (1986), 991-999.

D. C. Hoaglin, Discussion of "Graphical Perception: The Visual Decoding of Quantitative Information on Graphical Displays of Data" by W. S. Cleveland and R. McGill, Journal of the Royal Statistical Society, Series A, 150 (1987), 220.

B. Iglewicz and D. C. Hoaglin, "Use of Boxplots for Process Evaluation," Journal of Quality Technology, 19 (1987), 180-190.

D. C. Hoaglin and B. Iglewicz, "An Outlier-Labeling Procedure Based on the Biweight." Report AR-119, 10 December 1987.

M. Frigge, D. C. Hoaglin, and B. Iglewicz, "Some Implementations of the Boxplot," in R. M. Heiberger, ed., Computer Science and Statistics: Proceedings of the 19th Symposium on the

Interface. Alexandria, VA: American Statistical Association, 1987, pp. 296-300.

D. C. Hoaglin and B. Iglewicz, "Fine-Tuning Some Resistant Rules for Outlier Labeling," Journal of the American Statistical Association, 82 (1987), 1147-1149.

D. C. Hoaglin and B. Iglewicz, "Power of Some Outlier-Detection Rules under a Fixed-Outlier Model." Report AR-124, 27 May 1988.

D. C. Hoaglin and J. W. Tukey, "Empirical Bounds for Quantile-Based Estimates of g in the g -and- h Distributions." Report AR-125, 27 June 1988.

J. D. Emerson, D. C. Hoaglin, J. W. Tukey, and G. Y. Wong, "Exploring Some Adverb-Adjective Data," Chance, 1, 3 (Summer 1988), 42-48.

D. C. Hoaglin, "Using Leverage and Influence to Introduce Regression Diagnostics," College Mathematics Journal, 19 (1988), 387-401.

D. C. Hoaglin, "Transformations in Everyday Experience," Chance, 1, 4 (Fall 1988), to appear.

D. C. Hoaglin and J. H. Himes, Discussion of "Fitting Smoothed Centile Curves to Reference Data" by T. J. Cole, Journal of the Royal Statistical Society, Series A, 151 (1988), 412.

M. Frigge, D. C. Hoaglin, and B. Iglewicz, "Some Implementations of the Boxplot," The American Statistician, to appear.

J. H. Himes and D. C. Hoaglin, "Resistant Cross-Age Smoothing of Age-Specific Percentiles for Growth Reference Data," American Journal of Human Biology, to appear.

Participating Scientific Personnel

David C. Hoaglin (Research Associate)

Frederick Mosteller (Professor)

Cleo S. Youtz (Mathematical Assistant)

BIBLIOGRAPHY

- Frigge, Michael, Hoaglin, David C., and Iglewicz, Boris (1987). "Some Implementations of the Boxplot," in R. M. Heiberger, (Ed.), Computer Science and Statistics: Proceedings of the 19th Symposium on the Interface. Alexandria, VA: American Statistical Association, pp. 296-300.
- Frigge, Michael, Hoaglin, David C., and Iglewicz, Boris (1989). "Some Implementations of the Boxplot," The American Statistician, to appear.
- Himes, John H. and Hoaglin, David C. (1989). "Resistant Cross-Age Smoothing of Age-Specific Percentiles for Growth Reference Data," American Journal of Human Biology, to appear.
- Hoaglin, David C. (1985). "Summarizing Shape Numerically: The g -and- h Distributions," in David C. Hoaglin, Frederick Mosteller, and John W. Tukey (Eds.), Exploring Data Tables, Trends, and Shapes. New York: John Wiley & Sons, pp. 461-513.
- Hoaglin, David C. (1988a). "Using Leverage and Influence to Introduce Regression Diagnostics," College Mathematics Journal, 19, 387-401.
- Hoaglin, David C. (1988b). "Transformations in Everyday Experience," Chance, 1, 4 (Fall 1988), to appear.
- Hoaglin, David C. and Iglewicz, Boris (1987a). An Outlier-Labeling Procedure Based on the Biweight. Memorandum AR-119, Department of Statistics, Harvard University.
- Hoaglin, David C. and Iglewicz, Boris (1987b). "Fine-Tuning Some Resistant Rules for Outlier Labeling," Journal of the American Statistical Association, 82, 1147-1149.
- Hoaglin, David C. and Iglewicz, Boris (1988). Power of Some Outlier-Detection Rules under a Fixed-Outlier-Model. Memorandum AR-124, Department of Statistics, Harvard University.
- Hoaglin, David C., Iglewicz, Boris, and Tukey, John W. (1986). "Performance of Some Resistant Rules for Outlier Labeling," Journal of the American Statistical Association, 81, 991-999.
- Hoaglin, David C. and Kempthorne, Peter J. (1986). Discussion of "Influential Observations, High

- Leverage Points and Outliers in Linear Regression" by S. Chatterjee and A. S. Hadi, Statistical Science, 1, 408-412.
- Hoaglin, David C. and Tukey, John W. (1988). "Empirical Bounds for Quantile-Based Estimates of g in the g -and- h Distributions." Memorandum AR-125, Department of Statistics, Harvard University.
- Iglewicz, Boris and Hoaglin, David C. (1987). "Use of Boxplots for Process Evaluation," Journal of Quality Technology, 19, 180-190.
- Kafadar, Karen (1983). "The Efficiency of the Biweight as a Robust Estimator of Location," Journal of Research of the National Bureau of Standards, 88, 105-116.
- Lax, David A. (1985). "Robust Estimators of Scale: Finite-Sample Performance in Long-Tailed Symmetric Distributions," Journal of the American Statistical Association, 80, 736-741.
- Rocke, David M. (1983). "Robust Statistical Analysis of Interlaboratory Studies," Biometrika, 70, 421-431.
- Smith, Richard L. (1987). "Estimating Tails of Probability Distributions," Annals of Statistics, 15, 1174-1207.